

- 58 -

**CLAIMS:**

1. A method of computer data analysis using neural networks, the method including:  
generating a data representation using a data set, the data set including  
5 a plurality of attributes, wherein generating the data representation includes:  
modifying the data set using a training algorithm, wherein the  
training algorithm includes growing the data set; and  
performing convergence testing, wherein convergence testing  
checks for convergence of the training algorithm, and wherein the  
10 modifying of the data set is repeated until convergence of the training  
algorithm occurs; and  
displaying one or more subsets of the data set using the data  
representation.
2. A method according to claim 1, further including generating the  
15 data set using input data, and wherein generating the data set includes  
formatting the input data and initializing the formatted input data.
3. A method according to claim 2, wherein formatting the input data  
further includes creating a container class including a list of data vectors,  $\overline{\mathbf{D}}$ ,  
where  $\mathbf{d}_i$  is the  $i$ th vector in  $\overline{\mathbf{D}}$ , and  $\mathbf{d}_{i,j}$  is the  $j$ th element of vector  $i$ .
- 20 4. A method according to claim 2, wherein formatting the input data  
further includes data scaling and binarisation of at least a portion of the data  
set.
5. A method according to claim 4, wherein data scaling includes  
replacing each element in each data vector in the data set by a scaled  
25 representation of itself, where:  

$$\forall i \in [1, \text{card}(\mathbf{d})], \forall \mathbf{d}_i \in \overline{\mathbf{D}}$$

$$d_{i,j} = (d_{i,j} - i_{\min}) / (i_{\max} - i_{\min}).$$
6. A method according to claim 4 or 5, wherein binarisation includes  
converting attributes into one or more toggled attribute values.
- 30 7. A method according to any one of the preceding claims, wherein  
performing convergence testing includes testing condition  $q(t) < Q_c$ .

- 59 -

8. A method according to claim any one of claims 2 to 7, wherein initializing the formatted input data includes:

calculating an autocorrelation matrix,  $\mathbf{K}$  over the input data set  $\overline{\mathbf{D}}$ ,

$$\text{where } \mathbf{K} = \frac{1}{\text{card}(\overline{\mathbf{D}})} \sum_{\mathbf{d} \in \overline{\mathbf{D}}} \mathbf{d} \cdot \mathbf{d}^T ;$$

5 finding two longest eigenvectors of  $\mathbf{K}$ ,  $\mathbf{e}_1$  and  $\mathbf{e}_2$ , where  $|\mathbf{e}_1| > |\mathbf{e}_2|$ ; and  
initializing vector values of each element of the data set  $F$  by spanning it with element values of the eigenvectors.

9. A method according to claim 8, wherein initializing the vector values includes:

$$10 \quad F_{\langle F_R, 0 \rangle} := 0$$

$$F_{\langle F_R, F_C \rangle} := \mathbf{e}_1$$

$$F_{\langle 1, F_C \rangle} := \mathbf{e}_1 + \mathbf{e}_2$$

$$F_{\langle F_R, F_C \rangle} := \mathbf{e}_2$$

$$\forall c \in [2, F_C - 1], F_{\langle 1, c \rangle} := \frac{F_C}{F_C - c} F_{\langle 1, F_C \rangle} + \frac{F_C - c}{F_C} F_{\langle 1, 1 \rangle}$$

$$15 \quad \forall c \in [2, F_C - 1], F_{\langle R_C, c \rangle} := \frac{c}{F_C} F_{\langle F_R, F_C \rangle} + \frac{F_C - c}{F_C} F_{\langle F_R, 1 \rangle}$$

$$\forall r \in [2, F_R - 1], F_{\langle r, 1 \rangle} := \frac{c}{F_R} F_{\langle F_R, 1 \rangle} + \frac{F_R - r}{F_R} F_{\langle 1, 1 \rangle}$$

$$\forall r \in [2, F_R - 1], F_{\langle r, F_C \rangle} := \frac{r}{F_R} F_{\langle F_R, F_C \rangle} + \frac{F_R - r}{F_R} F_{\langle r, F_R \rangle}$$

$$\forall r \in [2, F_R - 1], \forall c \in [2, F_C - 1], F_{\langle r, c \rangle} := \frac{c}{F_C} F_{\langle r, F_C \rangle} + \frac{F_C - c}{F_C} F_{\langle r, 1 \rangle}.$$

10. A method according to any one of the preceding claims, wherein  
20 the data set includes a plurality of data set nodes, and wherein growing the data set includes:

finding  $K_q$  for each of the data set nodes, where  $K_q$  is the node

with the highest average quantization error,  $\arg \max_q \{ \overline{q}(t)_{K_q} \}$  for

- 60 -

each of the data set nodes, where  $\bar{q}(t)_{K_q} = \frac{1}{t-1} \sum_{i=1}^{t-1} q(i)_{K_q}$  is the average quantization error for node  $q$ , where:

$$K_x = \arg \max_x \{ \|K_q - K_{\langle r(q)-1, c(q) \rangle} \|, \|K_q - K_{\langle r(q)+1, c(q) \rangle} \| \}$$

$$K_y = \arg \max_y \{ \|K_q - K_{\langle r(q), c(q)-1 \rangle} \|, \|K_q - K_{\langle r(q), c(q)+1 \rangle} \| \}$$

5 if  $\|K_y - K_c\| < \|K_x - K_c\|$  then

$n_r = r(y)$  if  $r(y) < r(c)$ , else  $n_r = r(c)$ ; and

$n_c = c(y)$ ;

else  $n_r = r(y)$ ;  $n_c = c(x)$  if  $c(x) < c(c)$ , else  $n_c = c(c)$ ;

inserting a new row and column after row  $n_r$  and column  $n_c$ ; and

10 interpolating new attribute values for the newly inserted node

vectors using:  $K_{\langle r, n_c \rangle} = (K_{\langle r, n_c-1 \rangle} + K_{\langle r, n_c+1 \rangle}) \frac{\alpha}{2}$  and

$K_{\langle n_r, c \rangle} = (K_{\langle n_r-1, c \rangle} + K_{\langle n_r+1, c \rangle}) \frac{\alpha}{2}$ , where  $\alpha \in U(0,1)$ .

11. A method according to any one of the preceding claims, wherein the training algorithm further includes:

15  $t = t + 1$ ;

$\forall d \in \overline{D}$ ;

if ( $t < 50$  or *afterGrow*)

$$\wp_d = \arg \min_{\langle r, c \rangle, \forall r \in [1, P_R], \forall c \in [1, F_C]} \|d - F_{\langle r, c \rangle}\|_\rho$$

*afterGrow* = *false*

20 else

$\wp_d = \text{FindSCWS}(d)$

call function: *FindNeighborhoodPatterns*( $\wp$ )

call function: *BatchUpdateMatchVectors*

$$q(t) = \frac{1}{\text{card}(\overline{D})} \sum_i \|d - F_{\wp_d}\|_\rho$$
; and

25 if (*MayGrow*( $t$ ) and  $t < t_{\max}$ ), call function: *GrowKF*.

12. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes using a composite view

- 61 -

to view multiple attributes, wherein an additional attribute image is created, the additional attribute image displaying a union of a selected set of attributes.

13. A method according to claim 12, wherein using a composite view further includes:

5       constructing an attribute matrix; and  
      selecting a highest value for each attribute value from the selected set of attributes.

14. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes uses a range filter to  
10   select regions on the data representation and filter out nodes based on defined value ranges .

15. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes using a zooming function, wherein the zooming function includes:

15       making a selection of nodes to form a base reference of interest;  
      defining a set of data records from a second data set;  
      matching the second data set to the data representation;  
      flagging all records that are linked to the matched region; and  
      generating a second data representation using the flagged records.

20       16. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes using visual scaling, changing the minimum and maximum values used to calculate a colour progression used to visualize at least one of the plurality of attributes, and re-interpolating the active colour ranges over the new valid range of attribute  
25   values.

17. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes using a labeling engine to:

30       linking attribute columns in an input file to attributes in the data representation;  
      selecting attributes from the input file to be used for labelling;  
      determining with which row and column each row in the input file is associated; and  
      placing labels on the data representation.

- 62 -

18. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes using an advanced search function to:

- read a set of data records from a data source;
- 5 match attribute columns from the set of data records to attributes in the data representation; and
- display a list of all records that are associated with nodes that are part of the active selection on the data representation.

19. A method according to any one of the preceding claims, wherein displaying one or more subsets of the data set includes using equal distance averaging (EDA), wherein equal distance averaging includes:

finding the node vector,  $\mathbf{n}$ , in the data representation that most closely represents the input data vector,  $\mathbf{d}$ :

$$\mathbf{n} = \arg \min_{K_{r,c}} \left\{ \left\| K_{r,c} - \mathbf{d} \right\|_p \right\}, \forall r \in [1, K_R], \forall c \in [1, K_C]; \text{ and}$$

- 15 replacing missing entries in  $\mathbf{d}$  with the corresponding entries from  $\mathbf{n}$ .

20. A method according to claim 19, wherein the equal distance averaging further includes:

building a list of the data representation nodes values,  $\overline{\mathbf{M}}$ , such that for each element  $\mathbf{m}$  of  $\overline{\mathbf{M}}$ ,  $\left\| \mathbf{m} - \mathbf{d} \right\|_p = 0$ ;

- 20 if  $\overline{\mathbf{M}}$  is empty, then replace each missing entry in  $\mathbf{d}$  with corresponding entries in  $\mathbf{n}$ ; and

If  $\overline{\mathbf{M}}$  is not empty, then replace each missing entry in  $\mathbf{d}$  with the average value of the corresponding position of all the elements in  $\overline{\mathbf{M}}$ .

21. A method according to any one of the preceding claims, wherein the data representation includes a knowledge filter.

22. A method of computer data analysis using neural networks, the method including:

generating a data set  $\overline{\mathbf{D}}$ , the data set including a plurality of attributes and a plurality of data set nodes;

- 30 initializing the data set, initializing the data set including:

calculating an autocorrelation matrix,  $\mathbf{K}$  over the input data set  $\overline{\mathbf{D}}$ ,

$$\text{where } \mathbf{K} = \frac{1}{\text{card}(\overline{\mathbf{D}})} \sum_{\mathbf{d} \in \overline{\mathbf{D}}} \mathbf{d} \cdot \mathbf{d}^T ;$$

- 63 -

finding two longest eigenvectors of  $K$ ,  $e_1$  and  $e_2$ , where  $|e_1| > |e_2|$ ;  
and

initializing vector values of each element of a data representation  
 $F$  by spanning it with element values of the eigenvectors;

5 generating a data representation using a training algorithm, wherein the  
training algorithm includes growing the data set, growing the data set including:

finding  $K_q$  for each of the data set nodes, where  $K_q$  is the node  
with the highest average quantization error,  $\arg \max_q \{ \bar{q}(t)_{K_q} \}$  for

each of the data set nodes, where  $\bar{q}(t)_{K_q} = \frac{1}{t-1} \sum_{i=1}^{t-1} q(t)_{K_q}$  is the average

10 quantization error for node  $q$ , where:

$$K_x = \arg \max_x \{ \|K_q - K_{\langle r(q)-1, c(q) \rangle} \|, \|K_q - K_{\langle r(q)+1, c(q) \rangle} \| \}$$

$$K_y = \arg \max_y \{ \|K_q - K_{\langle r(q), c(q)-1 \rangle} \|, \|K_q - K_{\langle r(q), c(q)+1 \rangle} \| \}$$

if  $\|K_y - K_c\| < \|K_x - K_c\|$  then

$n_r = r(y)$  if  $r(y) < r(c)$ , else  $n_r = r(c)$ ; and

15  $n_c = c(y)$ ;

else  $n_r = r(y)$ ;  $n_c = c(x)$  if  $c(x) < c(c)$ , else  $n_c = c(c)$ ;

inserting a new row and column after row  $n_r$  and column  $n_c$ ;

interpolate new attribute values for the newly inserted node

vectors using:  $K_{\langle r, n_c \rangle} = (K_{\langle r, n_c-1 \rangle} + K_{\langle r, n_c+1 \rangle}) \frac{\alpha}{2}$  and

20  $K_{\langle n_r, c \rangle} = (K_{\langle n_r-1, c \rangle} + K_{\langle n_r+1, c \rangle}) \frac{\alpha}{2}$ , where  $\alpha \in U(0,1)$ ;

performing convergence testing, wherein convergence testing checks for  
convergence of the training algorithm, and wherein the training algorithm is  
repeated until convergence of the training algorithm occurs; and

displaying one or more subsets of the data set using the data  
25 representation.

23. A method according to claim 22, wherein initializing the vector  
values further includes:

$$F_{\langle F_R, 0 \rangle} := 0$$

- 64 -

$$F_{\langle F_R, F_C \rangle} := e_1$$

$$F_{\langle 1, F_C \rangle} := e_1 + e_2$$

$$F_{\langle F_R, F_C \rangle} := e_2$$

$$\forall c \in [2, F_C - 1], F_{\langle 1, c \rangle} := \frac{F_C}{F_C - c} F_{\langle 1, F_C \rangle} + \frac{F_C - c}{F_C} F_{\langle 1, 1 \rangle}$$

$$5 \quad \forall c \in [2, F_C - 1], F_{\langle R_C, c \rangle} := \frac{c}{F_C} F_{\langle F_R, F_C \rangle} + \frac{F_C - c}{F_C} F_{\langle F_R, 1 \rangle}$$

$$\forall r \in [2, F_R - 1], F_{\langle r, 1 \rangle} := \frac{c}{F_R} F_{\langle F_R, 1 \rangle} + \frac{F_R - r}{F_R} F_{\langle 1, 1 \rangle}$$

$$\forall r \in [2, F_R - 1], F_{\langle r, F_C \rangle} := \frac{r}{F_R} F_{\langle F_R, F_C \rangle} + \frac{F_R - r}{F_R} F_{\langle r, F_R \rangle}$$

$$\forall r \in [2, F_R - 1], \forall c \in [2, F_C - 1], F_{\langle r, c \rangle} := \frac{c}{F_C} F_{\langle r, F_C \rangle} + \frac{F_C - c}{F_C} F_{\langle r, 1 \rangle}.$$

24. A method according to claim 22 or 23, wherein the training  
10 algorithm further includes:

$$t = t + 1;$$

$$\forall d \in \overline{\mathbf{D}}$$

if ( $t < 50$  or *afterGrow*)

$$\wp_d = \arg \min_{\langle r, c \rangle, \forall r \in [1, F_R], \forall c \in [1, F_C]} \|d - F_{\langle r, c \rangle}\|_\rho$$

15 *afterGrow* = *false*

else

$$\wp_d = \text{FindSCWS}(d)$$

call function: *FindNeighborhoodPatterns*( $\overline{\wp}$ )

call function: *BatchUpdateMatchVectors*

$$20 \quad q(t) = \frac{1}{\text{card}(\overline{\mathbf{D}})} \sum_i (\|d - F_{\wp_d}\|_\rho).$$

if (*MayGrow*( $t$ ) and  $t < t_{\max}$ ), call function: *GrowKF*.

25. A method according to claim 22, 23, or 24, wherein performing convergence testing includes testing condition  $q(t) < Q_e$ .

26. A method according to any one of claims 22 to 25, wherein  
25 displaying one or more subsets of the data set includes using a composite view

- 65 -

to view multiple attributes, wherein an additional attribute image is created, the additional attribute image displaying a union of a selected set of attributes.

27. A method according to claim 26, wherein using a composite view further includes:

- 5        constructing an attribute matrix; and  
      selecting a highest value for each attribute value from the selected set of attributes.

28. A method according to any one of claims 22 to 27, wherein displaying one or more subsets of the data set includes uses a range filter to  
10       select regions on the data representation and filter out nodes based on defined value ranges .

29. A method according to any one of claims 22 to 28, wherein displaying one or more subsets of the data set includes using a zooming function, wherein the zooming function includes:

- 15       making a selection of nodes to form a base reference of interest;  
      defining a set of data records from a second data set;  
      matching the second data set to the data representation;  
      flagging all records that are linked to the matched region; and  
      generating a second data representation using the flagged records.

20       30. A method according to any one of claims 22 to 29, wherein displaying one or more subsets of the data set includes using visual scaling, changing the minimum and maximum values used to calculate a colour progression used to visualize at least one of the plurality of attributes, and re-interpolating the active colour ranges over the new valid range of attribute  
25       values.

31. A method according to any one of claims 22 to 30, wherein displaying one or more subsets of the data set includes using a labeling engine to:

- 30       linking attribute columns in an input file to attributes in the data representation;  
      selecting attributes from the input file to be used for labelling;  
      determining with which row and column each row in the input file is associated; and  
      placing labels on the data representation.



- 66 -

32. A method according to any one of claims 22 to 31, wherein displaying one or more subsets of the data set includes using an advanced search function to:

- read a set of data records from a data source;
- 5 match attribute columns from the set of data records to attributes in the data representation; and
- display a list of all records that are associated with nodes that are part of the active selection on the data representation.

33. A method according to any one of claims 22 to 32, wherein displaying one or more subsets of the data set includes using equal distance averaging (EDA), wherein equal distance averaging includes:

finding the node vector,  $\mathbf{n}$ , in the data representation that most closely represents the input data vector,  $\mathbf{d}$ :

$$\mathbf{n} = \arg \min_{K_{r,c}} \left\{ \left\| K_{r,c} - \mathbf{d} \right\|_p \right\}, \forall r \in [1, K_R], \forall c \in [1, K_C]; \text{ and}$$

- 15 replacing missing entries in  $\mathbf{d}$  with the corresponding entries from  $\mathbf{n}$ .

34. A method according to claim 33, wherein the equal distance averaging further includes:

building a list of the data representation nodes values,  $\overline{\mathbf{M}}$ , such that for each element  $\mathbf{m}$  of  $\overline{\mathbf{M}}$ ,  $\left\| \mathbf{m} - \mathbf{d} \right\|_p = 0$ ;

- 20 if  $\overline{\mathbf{M}}$  is empty, then replace each missing entry in  $\mathbf{d}$  with corresponding entries in  $\mathbf{n}$ ; and

If  $\overline{\mathbf{M}}$  is not empty, then replace each missing entry in  $\mathbf{d}$  with the average value of the corresponding position of all the elements in  $\overline{\mathbf{M}}$ .

35. A method according to any one of claims 22 to 34, wherein the data representation is a knowledge filter.

36. A method according to any one of the preceding claims wherein the data representation includes a latent model of the data set.

37. A system for performing data analysis using neural networks, the system including:

- 30 one or more processors;
- one or more memories coupled to the one or more processors; and

- 67 -

program instructions stored in the one or more memories, the one or more processors being operable to execute the program instructions, the program instructions including:

- generating a data representation using a data set, the data set including  
 5 a plurality of attributes, wherein generating the data representation includes:  
     modifying the data set using a training algorithm, wherein the training algorithm includes growing the data set; and  
     performing convergence testing, wherein convergence testing checks for convergence of the training algorithm, and wherein the  
 10 modifying of the data set is repeated until convergence of the training algorithm occurs; and  
     displaying one or more subsets of the data set using the data representation.

38. A system according to claim 37, wherein performing convergence  
 15 testing includes testing condition  $q(t) < Q_c$ .

39. A system according to claim 37 or 38, wherein the data set includes a plurality of data set nodes, and wherein growing the data set includes:

- finding  $K_q$  for each of the data set nodes, where  $K_q$  is the node  
 20 with the highest average quantization error,  $\arg \max_q \{ \bar{q}(t)_{K_q} \}$  for each of the data set nodes, where  $\bar{q}(t)_{K_q} = \frac{1}{t-1} \sum_{i=1}^{t-1} q(t)_{K_q}$  is the average quantization error for node  $q$ , where:

$$K_x = \arg \max_x \{ \|K_q - K_{\langle r(q)-1, c(q) \rangle} \|, \|K_q - K_{\langle r(q)+1, c(q) \rangle} \| \}$$

$$K_y = \arg \max_y \{ \|K_q - K_{\langle r(q), c(q)-1 \rangle} \|, \|K_q - K_{\langle r(q), c(q)+1 \rangle} \| \}$$

25 if  $\|K_y - K_c\| < \|K_x - K_c\|$  then

$n_r = r(y)$  if  $r(y) < r(c)$ , else  $n_r = r(c)$ ; and

$n_c = c(y)$ ;

else  $n_r = r(y)$ ;  $n_c = c(x)$  if  $c(x) < c(c)$ , else  $n_c = c(c)$ ;

inserting a new row and column after row  $n_r$  and column  $n_c$ ; and

- 68 -

interpolating new attribute values for the newly inserted node vectors using:  $K_{\langle r, n_r \rangle} = \left( K_{\langle r, n_r - 1 \rangle} + K_{\langle r, n_r + 1 \rangle} \right) \frac{\alpha}{2}$  and

$$K_{\langle n_r, c \rangle} = \left( K_{\langle n_r - 1, c \rangle} + K_{\langle n_r + 1, c \rangle} \right) \frac{\alpha}{2}, \text{ where } \alpha \in U(0,1).$$

40. A system according to claim 37, 38, or 39, wherein the training  
5 algorithm further includes:

$t = t + 1;$

$\forall d \in \overline{D};$

if ( $t < 50$  or *afterGrow*)

$$\phi_d = \arg \min_{\langle r, c \rangle, \forall r \in [1, F_R], \forall c \in [1, F_C]} \|d - F_{\langle r, c \rangle}\|_{\rho}$$

10 *afterGrow* = *false*

else

$\phi_d = \text{FindSCWS}(d)$

call function: *FindNeighborhoodPatterns*( $\overline{\phi}$ )

call function: *BatchUpdateMatchVectors*

15  $q(t) = \frac{1}{\text{card}(\overline{D})} \sum_i (\|d - F_{\phi_d}\|_{\rho}); \text{ and}$

if (*MayGrow*( $t$ ) and  $t < t_{\max}$ ), call function: *GrowKF*.

41. A system according to any one of claims 37 to 40, wherein the  
program instructions further include: displaying one or more subsets of the data  
set includes using a composite view to view multiple attributes, wherein an  
20 additional attribute image is created, the additional attribute image displaying a  
union of a selected set of attributes.

42. A system according to any one of claims 37 to 41, wherein the  
program instructions further include:

constructing an attribute matrix; and

25 selecting a highest value for each attribute value from the selected set of  
attributes.

43. A system according to any one of claims 37 to 42, wherein  
displaying one or more subsets of the data set includes uses a range filter to  
select regions on the data representation and filter out nodes based on defined  
30 value ranges.

- 69 -

44. A system according to any one of claims 37 to 43, wherein displaying one or more subsets of the data set includes using a zooming function, wherein the wherein the program instructions further include:

- 5 making a selection of nodes to form a base reference of interest;
- defining a set of data records from a second data set;
- matching the second data set to the data representation;
- flagging all records that are linked to the matched region; and
- generating a second data representation using the flagged records.

45. A system according to any one of claims 37 to 44, wherein displaying one or more subsets of the data set includes using visual scaling, wherein the wherein the program instructions further include:

- changing the minimum and maximum values used to calculate a colour progression used to visualize at least one of the plurality of attributes;
- and re-interpolating the active colour ranges over the new valid range of
- 15 attribute values.

46. A system according to any one of claims 37 to 45, wherein displaying one or more subsets of the data set includes using a labeling engine, wherein the program instructions further include:

- 20 linking attribute columns in an input file to attributes in the data representation;
- selecting attributes from the input file to be used for labelling;
- determining with which row and column each row in the input file is associated; and
- placing labels on the data representation.

47. A system according to any one of claims 37 to 46, wherein displaying one or more subsets of the data set includes using an advanced search engine, wherein the program instructions further include:

- displaying one or more subsets of the data set includes using an advanced search function to:
- 30 read a set of data records from a data source;
- match attribute columns from the set of data records to attributes in the data representation; and
- display a list of all records that are associated with nodes that are part of the active selection on the data representation.

- 70 -

48. A system according to any one of claims 37 to 47, wherein displaying one or more subsets of the data set includes using equal distance averaging (EDA), wherein the program instructions further include:

5 finding the node vector,  $\mathbf{n}$ , in the data representation that most closely represents the input data vector,  $\mathbf{d}$ :

$$\mathbf{n} = \arg \min_{K_{r,c}} \{ \|K_{r,c} - \mathbf{d}\|_p \}, \forall r \in [1, K_R], \forall c \in [1, K_C]; \text{ and}$$

replacing missing entries in  $\mathbf{d}$  with the corresponding entries from  $\mathbf{n}$ .

49. A system according to claim 48, wherein the program instructions further include:

10 building a list of the data representation nodes values,  $\overline{\mathbf{M}}$ , such that for each element  $\mathbf{m}$  of  $\overline{\mathbf{M}}$ ,  $\|\mathbf{m} - \mathbf{d}\|_p = 0$ ;

if  $\overline{\mathbf{M}}$  is empty, then replace each missing entry in  $\mathbf{d}$  with corresponding entries in  $\mathbf{n}$ ; and

15 If  $\overline{\mathbf{M}}$  is not empty, then replace each missing entry in  $\mathbf{d}$  with the average value of the corresponding position of all the elements in  $\overline{\mathbf{M}}$ .

50. A system according to any one of claims 37 to 49, wherein the data representation includes a knowledge filter.

51. A system according to any one of claims 37 to 50, wherein the data representation includes a latent model of the data set.

20 52. A computer program product for computer data analysis using neural networks, the computer program product including:

computer-readable program code for generating a data representation using a data set, the data set including a plurality of attributes, wherein generating the data representation includes:

25 modifying the data set using a training algorithm, wherein the training algorithm includes growing the data set; and

performing convergence testing, wherein convergence testing checks for convergence of the training algorithm, and wherein the modifying of the data set is repeated until convergence of the training algorithm occurs; and

30 computer-readable program code for displaying one or more subsets of the data set using the data representation.

53. A computer program product according to claim 52, wherein the data set includes a plurality of data set nodes, and the computer program product further including computer-readable program code for growing the data set including:

- 5 finding  $K_q$  for each of the data set nodes, where  $K_q$  is the node with the highest average quantization error,  $\arg \max_q \{ \bar{q}(t)_{K_q} \}$  for each of the data set nodes, where  $\bar{q}(t)_{K_q} = \frac{1}{t-1} \sum_{i=1}^{t-1} q(t)_{K_q}$  is the average quantization error for node  $q$ , where:

$$K_x = \arg \max_x \{ \|K_q - K_{\langle r(q)-1, c(q) \rangle} \|, \|K_q - K_{\langle r(q)+1, c(q) \rangle} \| \}$$

$$10 \quad K_y = \arg \max_y \{ \|K_q - K_{\langle r(q), c(q)-1 \rangle} \|, \|K_q - K_{\langle r(q), c(q)+1 \rangle} \| \}$$

if  $\|K_y - K_c\| < \|K_x - K_c\|$  then

$n_r = r(y)$  if  $r(y) < r(c)$ , else  $n_r = r(c)$ ; and

$n_c = c(y)$ ;

else  $n_r = r(y)$ ;  $n_c = c(x)$  if  $c(x) < c(c)$ , else  $n_c = c(c)$ ;

- 15 inserting a new row and column after row  $n_r$  and column  $n_c$ ; and interpolating new attribute values for the newly inserted node vectors using:  $K_{\langle r, n_c \rangle} = (K_{\langle r, n_c-1 \rangle} + K_{\langle r, n_c+1 \rangle}) \frac{\alpha}{2}$  and

$$K_{\langle n_r, c \rangle} = (K_{\langle n_r-1, c \rangle} + K_{\langle n_r+1, c \rangle}) \frac{\alpha}{2}, \text{ where } \alpha \in U(0,1).$$

54. A computer program product according to claim 52 or 53, wherein  
20 the data representation includes a latent model of the data set.

55. An apparatus for performing data analysis using neural networks, the apparatus including:

means for representing a data set, the data set including a plurality of attributes;

- 25 means for generating the representation means using the data set, wherein generating the representation means includes:

modifying the data set using a training algorithm, wherein the training algorithm includes growing the data set; and

- 72 -

performing convergence testing, wherein convergence testing checks for convergence of the training algorithm, and wherein the modifying of the data set is repeated until convergence of the training algorithm occurs; and

5 means for displaying one or more subsets of the data set using the modified data representation.

56. An apparatus according to claim 55, wherein the representation means includes a latent model of the data set.

57. A method of computer data analysis using neural networks  
10 substantially as herein described with reference to the accompanying drawings.

58. A system for performing data analysis using neural networks substantially as herein described with reference to the accompanying drawings.

59. A computer program product for computer data analysis using neural networks substantially as herein described with reference to the  
15 accompanying drawings.

60. An apparatus for performing data analysis using neural networks substantially as herein described with reference to the accompanying drawings.